

Prednáška 5 - Jednoduché pravidlá, Chomského normálny tvar, CYK algoritmus

Ing. Viliam Hromada, PhD.

C-510
Ústav informatiky a matematiky
FEI STU

`viliam.hromada@stuba.sk`



Transformácie bezkontextových gramatík

- Minulý týždeň sme sa začali baviť o transformáciách gramatík
- Spomenuli sme úpravu odstránenia nadbytočných symbolov gramatiky, t.j. jej úpravu do tzv. redukovaného tvaru.
- Rovnako sme spomínali úpravu odstránenia ε -pravidiel
- Dnes si uvedieme ďalšiu úpravu - odstránenie jednoduchých pravidiel.
- Ukážeme si úpravu gramatiky do Chomského normálnej formy.
- A jej praktické využitie v CYK algoritme.



Motivácia úprav gramatiky

- Nech je daná bezkontextová gramatika $G = (N, T, P, S)$ a reťazec $w \in T^*$ zložený z terminálnych symbolov.
- Uvažujme rozhodovací problém: $w \in L(G)$?. Teda, má reťazec w v gramatike G deriváciu?
- Existujú algoritmy, ktoré dokážu riešiť tento problém s **polynomiálnou zložitosťou** (pre bezkontextové a regulárne gramatiky).
- Pre ich aplikáciu je však často potrebné upraviť príslušnú gramatiku G na ekvivalentnú gramatiku spĺňajúcu isté ďalšie vlastnosti - a to vieme docieľiť pomocou spomínaných úprav (ako napr. odstránenie ε -pravidiel alebo odstránenie **jednoduchých pravidiel**).



Odstránenie jednoduchých pravidiel a cyklov

Definícia

*Nech $G = (N, T, P, S)$ je bezkontextová gramatika. Potom pravidlá tvaru $A \rightarrow B$, kde $A, B \in N$ sa nazývajú **jednoduché pravidlá**.*



Odstránenie jednoduchých pravidiel

- Každé jednoduché pravidlo typu $A \rightarrow B$ odstránime a zároveň do gramatiky doplníme ďalšie pravidlá, aby sa nezmenil výsledný jazyk generovaný gramatikou.
- Zjednodušené, ak sa z A dalo odvodiť B , tak po odstránení pravidla $A \rightarrow B$ musíme doplniť pravidlá, ktoré majú na ľavej strane A a na pravej strane všetky možné pravé strany pravidiel typu $B \rightarrow$



Príklad

Gramatiku s jednoduchými pravidlami

$$S' \rightarrow S \mid \varepsilon$$

$$S \rightarrow aSb \mid A \mid ab$$

$$A \rightarrow bAa \mid ba$$

Upravíme na:

$$S' \rightarrow aSb \mid bAa \mid ba \mid ab \mid \varepsilon$$

$$S \rightarrow aSb \mid bAa \mid ba \mid ab$$

$$A \rightarrow bAa \mid ba$$



Odstránenie jednoduchých pravidiel

Pre každý neterminál A vytvoríme množinu N_A neterminálov, ktoré sa z A dajú odvodiť, t.j.:

$$N_A = \{B \mid A \Rightarrow^* B, B \in N\} \quad (1)$$



Vytvorenie množiny N_A

Vstup: bezkontextová gramatika $G = (N, T, P, S)$ bez ε -pravidiel, neterminál A

Výstup: množina N_A (4)

1: $N_A \leftarrow \{A\}$

2: **opakuj**

3: $\dot{N}_A \leftarrow N_A;$

4: $N_A \leftarrow \dot{N}_A \cup \{C \mid B \rightarrow C \in P \wedge B \in \dot{N}_A\}$

5: **pokiaľ** $N_A \neq \dot{N}_A$



Odstránenie jednoduchých pravidiel

- Na základe množín N_A teraz môžeme odstrániť jednoduché pravidlá.
- Pre každý neterminál vieme, na aké rôzne iné neterminály sa dá prepísať.
- Ak chceme odstrániť jednoduché pravidlá, musíme znovu zabezpečiť, aby to nezmenilo výsledný jazyk, ktorý gramatika generuje.
- To docielime tak, že pre všetky $B \in N_A$ a všetky pravidlá $B \rightarrow \alpha$, ktoré nie sú jednoduché, pridáme do upravenej množiny pravidiel aj pravidlá $A \rightarrow \alpha$.



Odstránenie jednoduchých pravidiel - príklad

Nech je daná gramatika (bola výsledkom úpravy odstránenia ε -pravidiel):

$$\acute{S} \rightarrow S \mid \varepsilon$$

$$S \rightarrow aAbB \mid AC \mid abB \mid A \mid C$$

$$A \rightarrow Aa \mid a$$

$$B \rightarrow bBb \mid ab$$

$$C \rightarrow aCA \mid A \mid aC \mid aA \mid a$$

Pre jednotlivé neterminály sú množiny N_A (postupne):

$$\acute{S} : \{\acute{S}\}, \{\acute{S}, S\}, \{\acute{S}, S, A, C\} = N_{\acute{S}}$$

$$S : \{S\}, \{S, A, C\} = N_S$$

$$A : \{A\} = N_A$$

$$B : \{B\} = N_B$$

$$C : \{C\} = \{C, A\} = N_C$$



Odstránenie jednoduchých pravidiel - príklad

Po odstránení jednoduchých pravidiel dostávame pravidlá:

$$\hat{S} \rightarrow \varepsilon \mid aAbB \mid AC \mid abB \mid Aa \mid a \mid aCA \mid aC \mid aA$$

$$S \rightarrow aAbB \mid AC \mid abB \mid Aa \mid a \mid aCA \mid aC \mid aA$$

$$A \rightarrow Aa \mid a$$

$$B \rightarrow bBb \mid ab$$

$$C \rightarrow aCA \mid aC \mid aA \mid a \mid Aa$$

T.j. gramatiku bez jednoduchých pravidiel.



Praktické použitie úprav

Spomenuli sme teda 3 typy úprav bezkontextových gramatík:

- Odstránenie ε -pravidiel
- Odstránenie jednoduchých pravidiel
- Odstránenie nadbytočných a nedostupných symbolov.

Tieto 3 typy úprav v praxi používame napríklad pri tvorbe **Chomského normálneho tvaru gramatiky**.



Chomského normálny tvar (forma) gramatiky

Definícia

Bezkontextová gramatika $G = (N, T, P, S)$ je v **Chomského normálnom tvare**, ak každé pravidlo je v jednom z tvarov:

1. $A \rightarrow BC$ pre $A \in N; B, C \in (N - \{S\})$;
2. $A \rightarrow a$ pre $A \in N, a \in T$;
3. $S \rightarrow \varepsilon$.



Úprava do CHNT

Vstup: Bezkontextová gramatika $G = (N, T, P, S)$

Výstup: Gramatika G v CHNT

- 1: Odstráň z gramatiky ε -pravidlá.
- 2: Uprav gramatiku, aby počiatočný neterminál nebol v pravej strane žiadneho pravidla.
- 3: Odstráň z gramatiky jednoduché pravidlá.
- 4: Odstráň z gramatiky nadbytočné a nedostupné symboly.
- 5: **pre všetky** pravidlá tvaru $A \rightarrow X_1 X_2 \dots X_n, n \geq 3$ **rob**
- 6: nahraď pravidlami $A \rightarrow X_1 A_1, A_1 \rightarrow X_2 A_2, \dots, A_{n-2} \rightarrow X_{n-1} X_n$, kde A_1, \dots, A_{n-2} sú nové neterminály.
- 7: **koniec pre**
- 8: **pre všetky** pravidlá tvaru $A \rightarrow Y_1 Y_2$, kde $Y_1 \in T \vee Y_2 \in T$ **rob**
- 9: nahraď na pravej strane pravidiel terminály $Y_1 (Y_2)$ novými neterminálmi $V_1 (V_2)$ a pridaj pravidlo (pravidlá) $V_1 \rightarrow Y_1$, resp. $V_2 \rightarrow Y_2$
- 10: **koniec pre**



Krok č. 2 z algoritmu

V algoritme sa po odstránení ε -pravidiel robí úprava: *Uprav gramatiku, aby počiatočný neterminál nebol v pravej strane žiadneho pravidla..*

- V prípade, že sa pri odstránení ε pridal do gramatiky nový neterminál \acute{S} a pravidlá $\acute{S} \rightarrow \varepsilon \mid S$, tak tento krok je irelevantný.
- V prípade, že sa v gramatike jej **aktuálny** počiatočný neterminál S nachádza v niektorej pravej strane niektorého pravidla, **pridá sa nový počiatočný neterminál \acute{S}** a pravidlo $\acute{S} \rightarrow S$.



CHNT - príklad

Gramatika $G = (\{S, A\}, \{a, b, c\}, P, S)$ s pravidlami:

$$S \rightarrow aASa \mid c$$

$$A \rightarrow Abc \mid c$$

- 1) ε -pravidlá nie je potrebné odstrániť - žiadne neobsahuje.
- 2) Zavedieme symbol \hat{S} a pridáme pravidlo:

$$\hat{S} \rightarrow S$$

$$S \rightarrow aASa \mid c$$

$$A \rightarrow Abc \mid c$$



CHNT - príklad (pokr.)

3,4) Po odstránení jednoduchých pravidiel a nadbytočných symbolov:

$$\acute{S} \rightarrow aASa \mid c$$

$$S \rightarrow aASa \mid c$$

$$A \rightarrow Abc \mid c$$

5-7) Úprava pravidiel dlhších ako 3 symboly:

$$\acute{S} \rightarrow aA_1 \mid c$$

$$S \rightarrow aA_1 \mid c$$

$$A \rightarrow AA_3 \mid c$$

$$A_1 \rightarrow AA_2$$

$$A_2 \rightarrow Sa$$

$$A_3 \rightarrow bc$$



CHNT - príklad (pokr.)

8-10) Odstránenie terminálov z pravidiel s pravou stranou dĺžky 2

$$\dot{S} \rightarrow V_1 A_1 \mid c$$

$$S \rightarrow V_1 A_1 \mid c$$

$$A \rightarrow AA_3 \mid c$$

$$A_1 \rightarrow AA_2$$

$$A_2 \rightarrow SV_1$$

$$A_3 \rightarrow V_2 V_3$$

$$V_1 \rightarrow a$$

$$V_2 \rightarrow b$$

$$V_3 \rightarrow c$$



CHNT - poznámky

- Ak vieme v danej gramatike odvodiť slovo w , potom dĺžka jeho odvodu v CHNT je $2|w| - 1$.
- Strom odvodu (derivačný strom) slova v gramatike v CHNT je binárny strom.



Praktické využitie Chomského normálneho tvaru

- Chomského normálny tvar sa v praxi využíva v tzv. CYK algoritme.
- CYK algoritmus je **polynomiálnym** algoritmom na zistenie, či **ľubovoľná bezkontextová gramatika** $G = (N, T, P, S)$ generuje zadaný **reťazec terminálov** $w \in T^*$, teda či $w \in L(G)$.



CYK algoritmus

- Predpokladajme, že máme danú bezkontextovú gramatiku $G = (N, T, P, S)$ v **Chomského normálnom tvare**.
- Chceme zistiť, či dané slovo $w = a_1 \dots a_n \in L(G)$, $w \in T^*$.
- Algoritmus CYK (Cocke-Younger-Kasami) to vie zistiť so zložitou $O(n^3 \cdot |P|)$, kde n je dĺžka analyzovaného slova pre **hocijakú** bezkontextovú gramatiku a $|P|$ počet pravidiel gramatiky.
- Základnou ideou algoritmu je konštrukcia množín:

$$N_{i,j} = \{A \mid A \Rightarrow^* a_i a_{i+1} \dots a_j, A \in N\},$$

$$1 \leq i \leq j \leq n.$$

- Ak $S \in N_{1,n}$ potom slovo $w \in L(G)$.



CYK algoritmus

- Postupne sa hľadajú neterminály, ktoré vedia generovať podreťazce daného slova w .
- Na začiatku sa určia neterminály generujúce podreťazce dĺžky 1.
 - Keďže gramatika je v CNF, dané neterminály sa určia priamo z pravidiel.
- Následne sa využíva nasledovná vlastnosť:
 - $B \in N_{j,k}$, t.j. $B \Rightarrow^* a_j \dots a_k$
 - $C \in N_{k+1,j+i}$, t.j. $C \Rightarrow^* a_{k+1} \dots a_{j+i}$,
a existuje pravidlo $A \rightarrow BC$, potom $A \in N_{j,j+i}$, pretože
 $A \Rightarrow BC \Rightarrow^* a_j \dots a_k a_{k+1} \dots a_{j+i}$.



CYK algoritmus I

Vstup: Bezkontextová gramatika $G = (N, T, P, S)$ v CNF, slovo $w \neq \varepsilon$

Výstup: Zistenie, či $w \in L(G)$

- 1: $w = a_1 \dots a_n$
- 2: **pre všetky** $i \in \{1, \dots, n\}$ **rob**
- 3: **pre všetky** $j \in \{i, \dots, n\}$ **rob**
- 4: $N_{i,j} \leftarrow \emptyset$
- 5: **koniec pre**
- 6: **koniec pre**
- 7: **pre všetky** $i \in \{1, \dots, n\}$ **rob**
- 8: **pre všetky** pravidlá tvaru $A \rightarrow a_i$ **rob**
- 9: $N_{i,i} \leftarrow N_{i,i} \cup \{A\}$
- 10: **koniec pre**



CYK algoritmus II

11: **koniec pre**

12: **pre všetky** $i \in \{1, 2, \dots, n - 1\}$ **rob**

13: **pre všetky** $j \in \{1, 2, \dots, n - i\}$ **rob**

14: **pre všetky** $k \in \{j, \dots, j + i - 1\}$ **rob**

15: **pre všetky** pravidlá $A \rightarrow BC$, kde $B \in N_{j,k}$ a zároveň $C \in N_{k+1,j+i}$ **rob**

16: $N_{j,j+i} \leftarrow N_{j,j+i} \cup \{A\}$

17: **koniec pre**

18: **koniec pre**

19: **koniec pre**

20: **koniec pre**

21: **ak** $S \in N_{1,n}$ **potom**

22: **vrát'** „Slovo patrí do jazyka. “



CYK algoritmus III

23: **inak**

24: **vrát'** „Slovo nepatrí do jazyka. “

25: **koniec ak**



CYK algoritmus - príklad

Je daná gramatika:

$$S \rightarrow AB \mid AC \mid a$$

$$A \rightarrow DA \mid a$$

$$B \rightarrow b \mid d$$

$$C \rightarrow BC \mid c$$

$$D \rightarrow d$$

Zistite pomocou CYK algoritmu, či slovo $dab \in L(G)$.



CYK algoritmus - príklad (pokr.)

Najprv sa zostroja množiny

$N_{1,1} = \{B, D\}$, kvôli pravidlám $B \rightarrow d$ a $D \rightarrow d$.

$N_{2,2} = \{S, A\}$, kvôli pravidlám $S \rightarrow a$ a $A \rightarrow a$.

$N_{3,3} = \{B\}$, kvôli pravidlu $B \rightarrow b$.

Ďalej podľa iterácii:

$i = 1, j = 1, k = 1$ $A \rightarrow DA, D \in N_{1,1}, A \in N_{2,2}$
teda $A \in N_{1,2}$

$i = 1, j = 2, k = 2$ $S \rightarrow AB, A \in N_{2,2}, B \in N_{3,3}$
teda $S \in N_{2,3}$

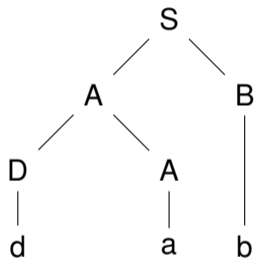
$i = 2, j = 1, k = 1$ $N_{1,1} = \{B, D\}, N_{2,3} = \{S\}$
do $N_{1,3}$ nepribudne nič

$i = 2, j = 1, k = 2$ $S \rightarrow AB, A \in N_{1,2}, B \in N_{3,3}$
teda $S \in N_{1,3}$

Keďže $S \in N_{1,3}$, tak $dab \in L(G)$.



Na základe príslušných množín vieme teoreticky aj nájsť derivačný strom, resp. odvodenie:



- Základnou úlohou CYK je zistiť, či slovo patrí do jazyka generovaného gramatikou, avšak v podstate sa dá použiť aj na zostrojenie derivačného stromu (na základe množín $N_{i,j}$).
- Výhodou CYK algoritmu je, že vie spracovať všetky bezkontextové gramatiky.
 - Aj nejednoznačné, aj jednoznačné
 - Aj nedeterministické, aj deterministické
- Nevýhodou CYK algoritmu je, že je potrebné, aby gramatika bola v CHNT a taktiež má kubickú zložitosť vzhľadom na dĺžku vstupného slova.
- Existujú aj ďalšie algoritmy, ktoré síce nevedia spracovať všetky bezkontextové gramatiky, avšak majú lineárnu zložitosť vzhľadom na dĺžku vstupného slova.



Použitá literatúra

Aho, A., Lam, M., Sethi, R., Ullman, J.: *Compilers: Principles, techniques and tools.*

Dedera, Ľ: *Počítačové jazyky a ich spracovanie.*

Linz, P.: *An Introduction to Formal Languages and Automata.*

